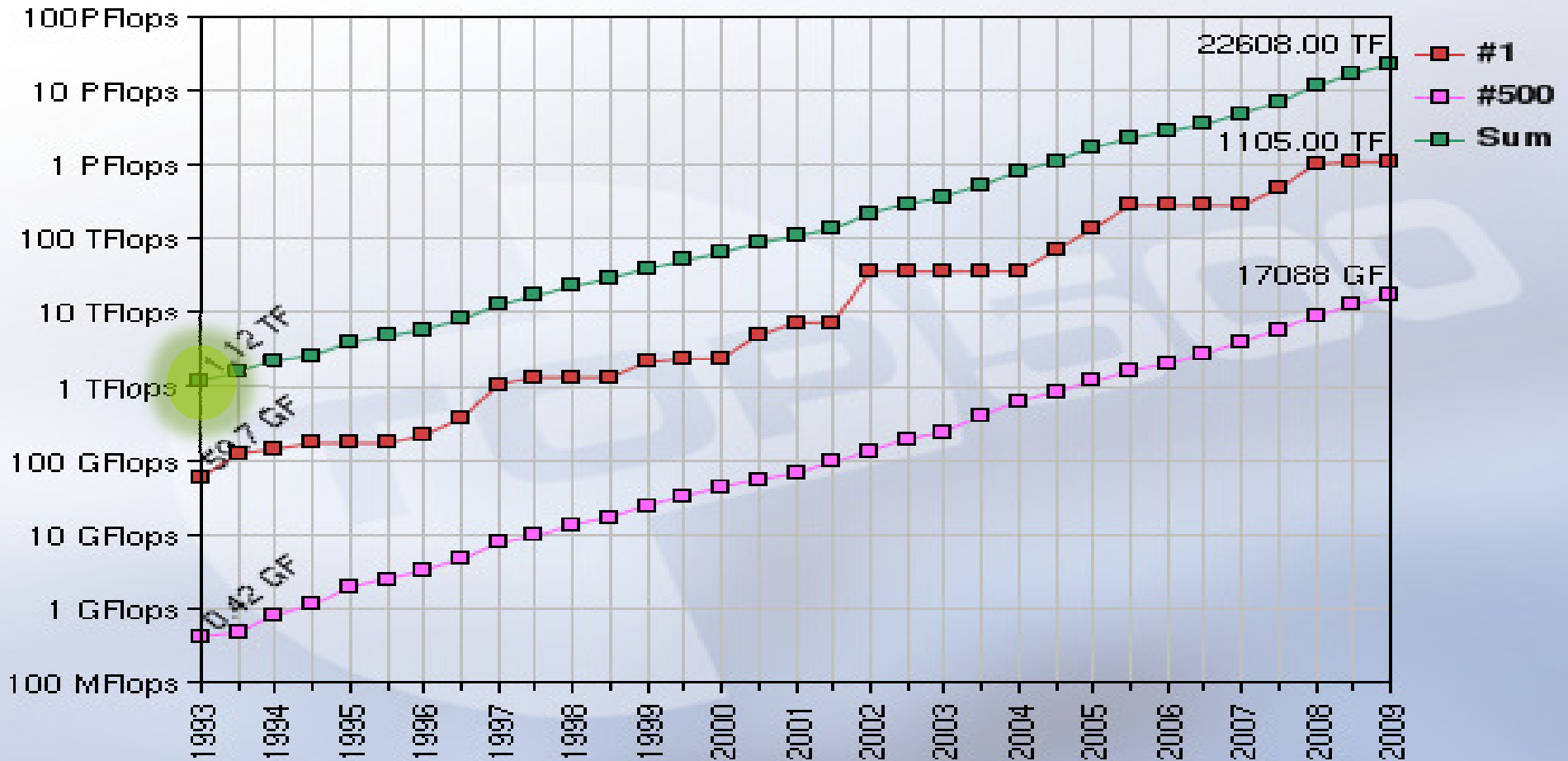# The Power of Heterogeneous Parallel Computing: Computational Graphics
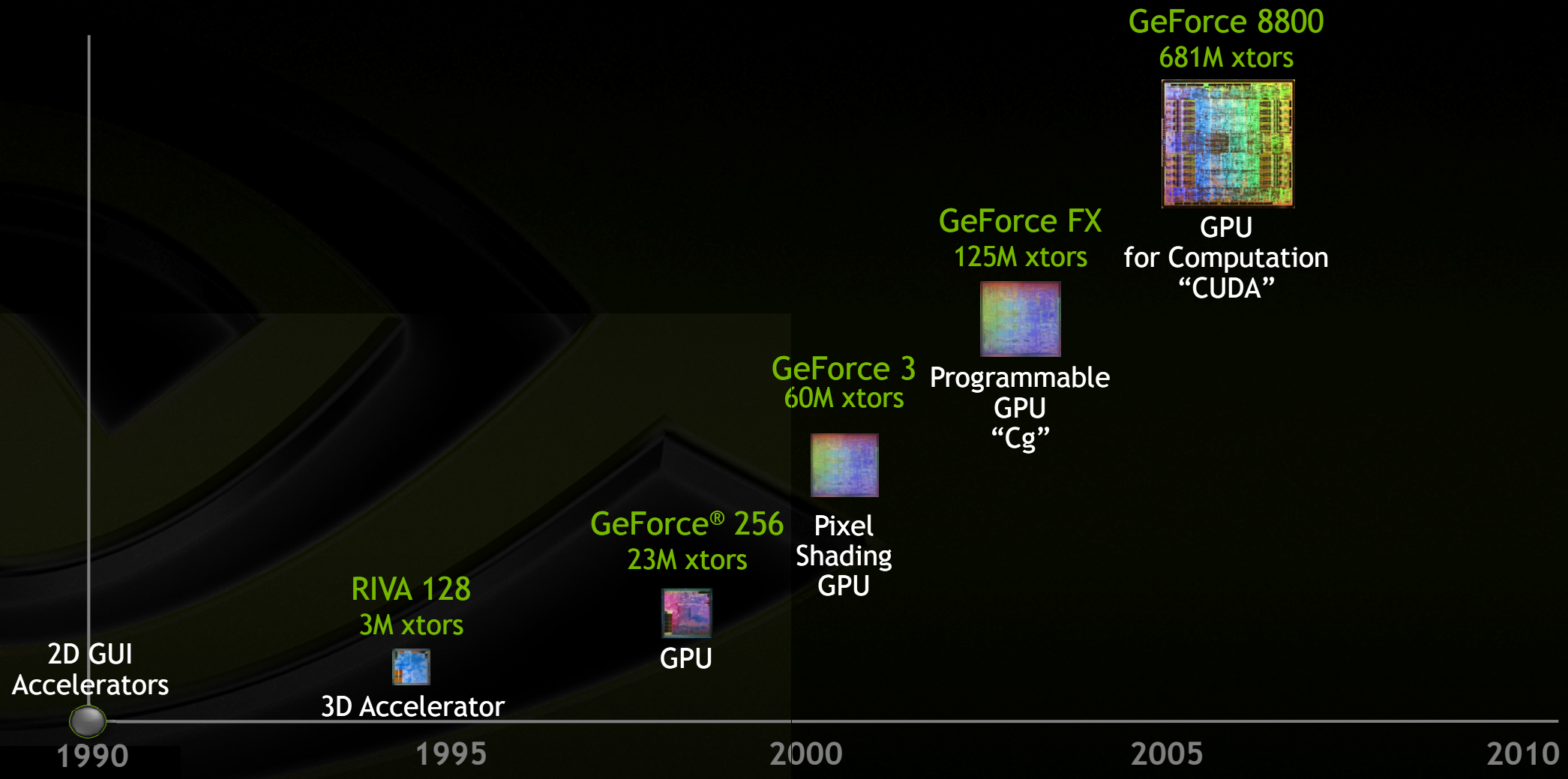
**David B. Kirk, NVIDIA Fellow**

# Performance Development
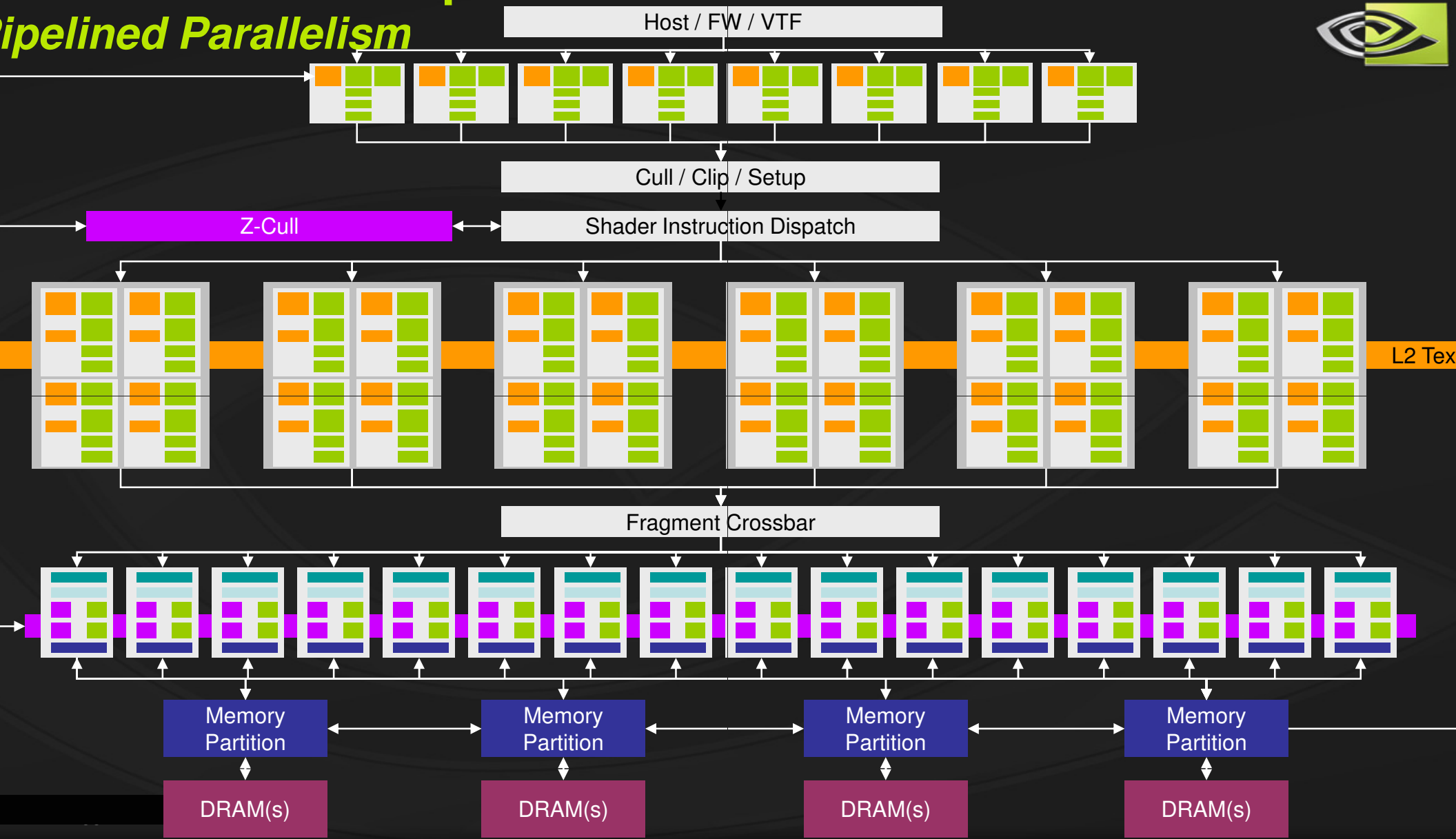
# NVIDIA Technology Evolution

**GeForce 8800**
681M xtors

GPU
for Computation
"CUDA"

**GeForce FX**
125M xtors

Programmable
GPU
"Cg"

**GeForce 3**
60M xtors

Pixel
Shading
GPU

**GeForce® 256**
23M xtors

GPU

**RIVA 128**
3M xtors

3D Accelerator

2D GUI
Accelerators

1990    1995    2000    2005    2010

**GeForce 7800 - Graphics**
*Pipelined Parallelism*

Host / FW / VTF

Cull / Clip / Setup

Z-Cull

Shader Instruction Dispatch

L2 Tex

Fragment Crossbar

Memory Partition

Memory Partition

Memory Partition

Memory Partition

DRAM(s)
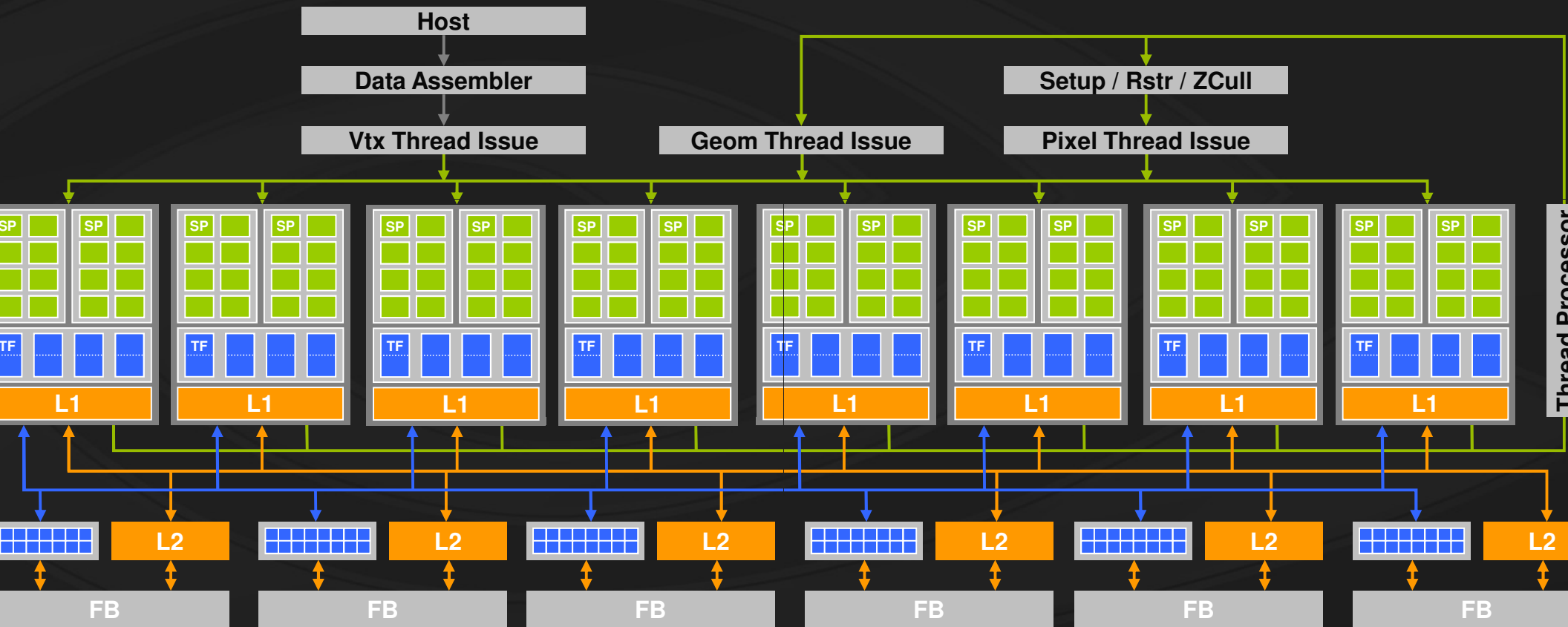
DRAM(s)

DRAM(s)

DRAM(s)

# GF8800 replaces the pipeline model

- The future of GPUs is programmable processing
- So – build the architecture around the processor(s)

# Revolutionizing Computer Graphics

# Computer Graphics – A Study in Parallelism

1 Frame Time
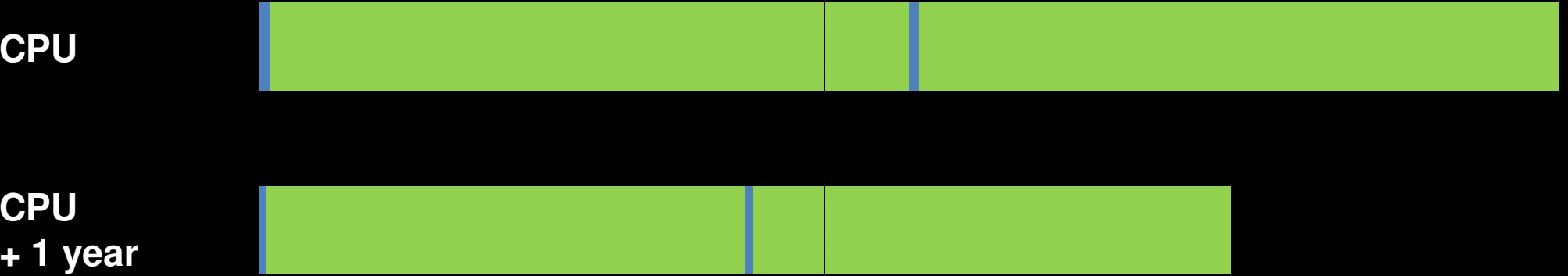
**CPU**

**Programmable Shading**

Nehalem CPU:
 3 Ghz
 4 cores
 4 way SIMD
 2 FLOPS/cycle
96 GFLOPS

2,300,000 pixels/frame

x 3 depth complexity

x 100 shader inst./component

x 1.5 FLOPS/inst.

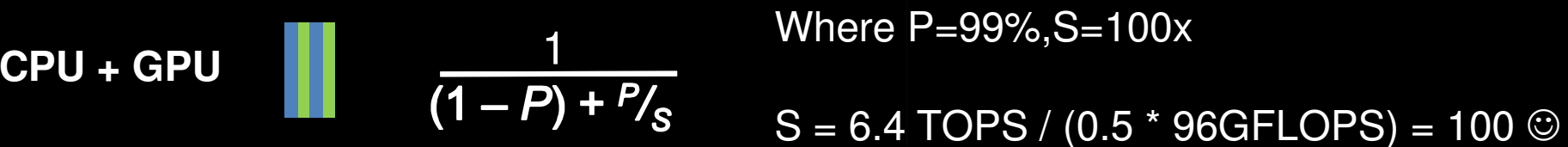x 4 components/pixel

x 60 frames/sec

x 2 stereoscopic

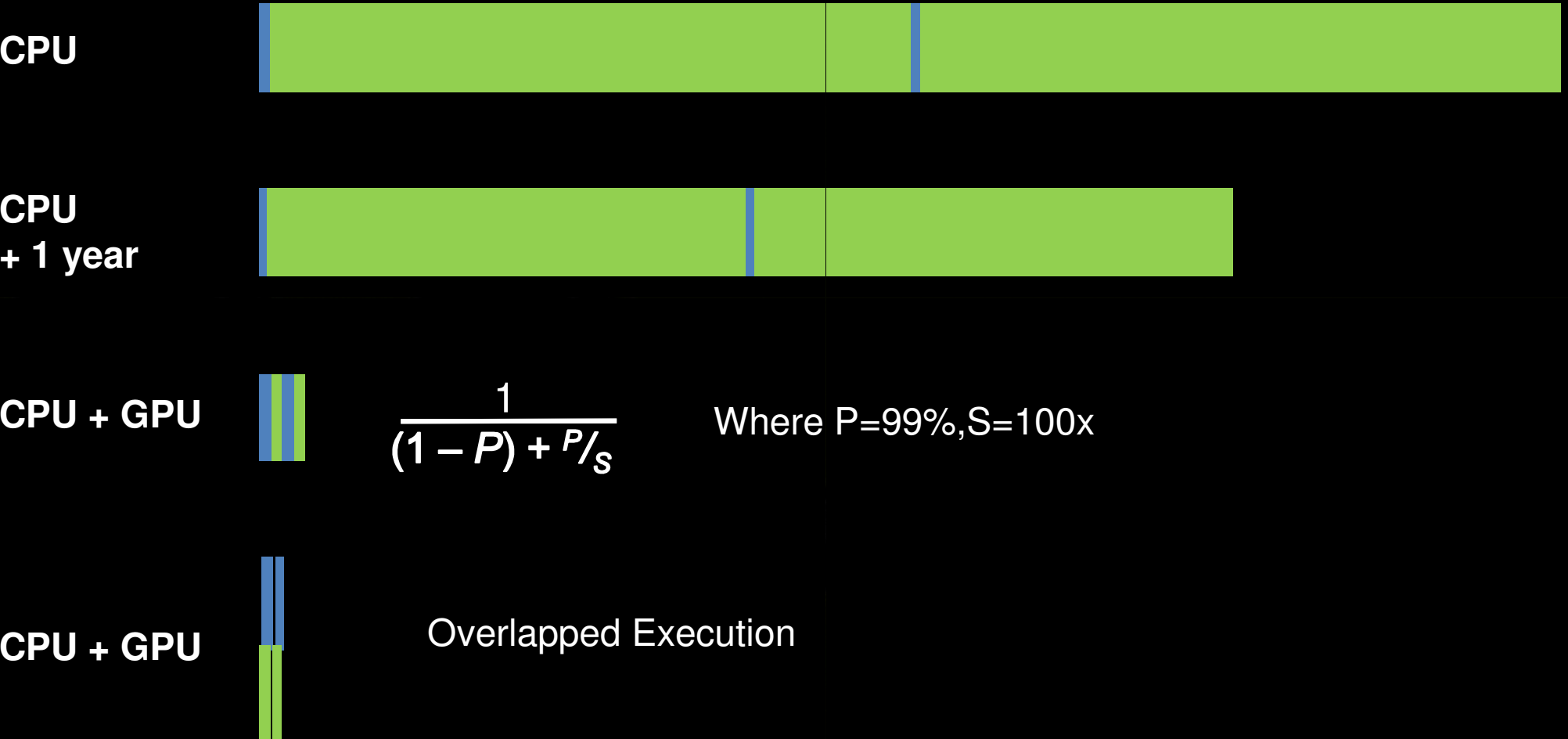= 500 shader GFLOPS
(approx. 10% of graphics ops)

# Computer Graphics – A Study in Parallelism

**CPU**

**CPU
+ 1 year**

# Computer Graphics – A Study in Parallelism

**CPU**

**CPU + 1 year**

**CPU + GPU**

$$\frac{1}{(1 - P) + {}^{P}\!/_{S}}$$

Where P=99%,S=100x

S = 6.4 TOPS / (0.5 * 96GFLOPS) = 100 ☺

# Computer Graphics – A Study in Parallelism

**CPU**

**CPU + 1 year**

**CPU + GPU**

$$\frac{1}{(1 - P) + {^P}\!/_S}$$

Where P=99%,S=100x
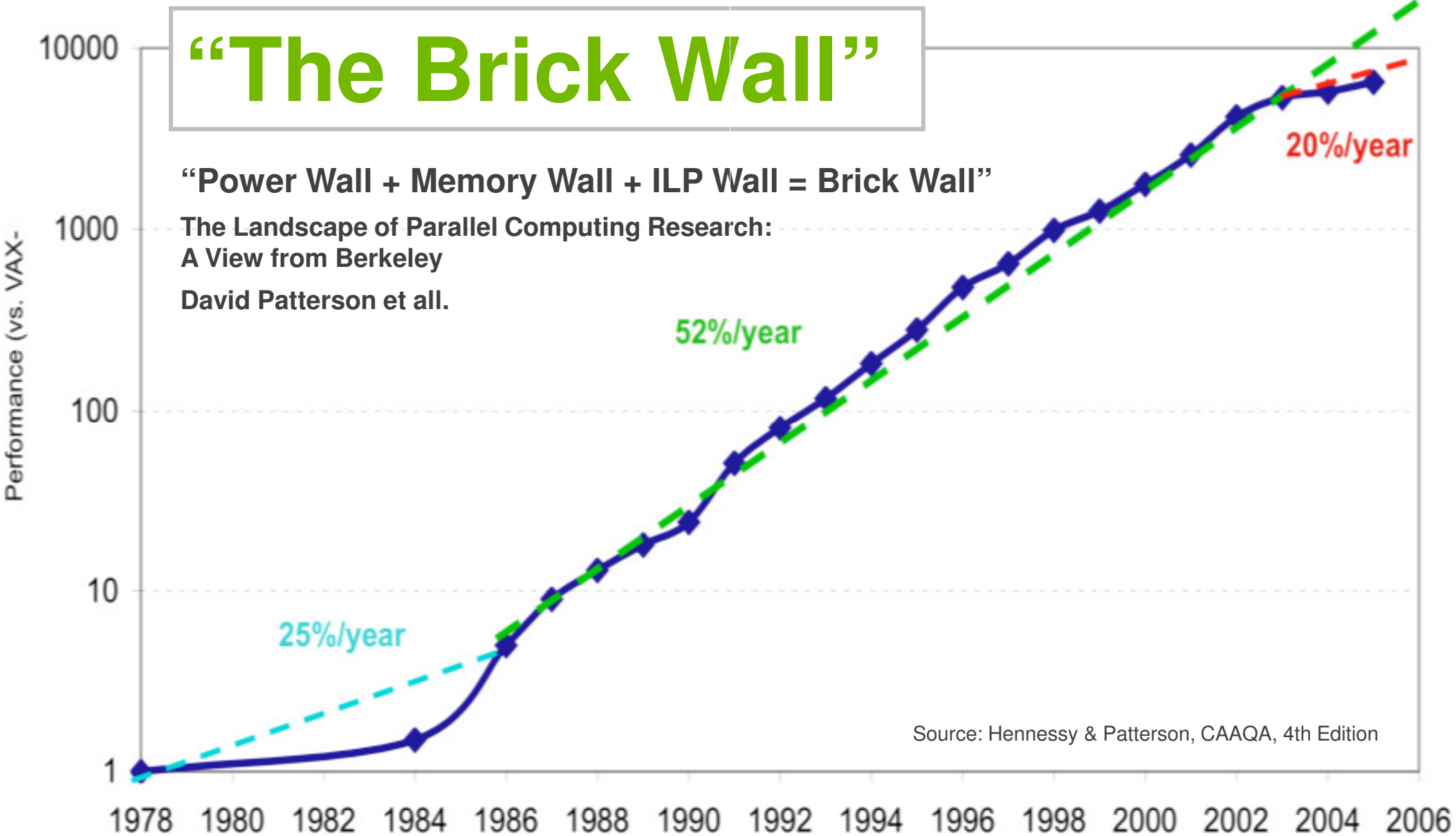
**CPU + GPU**

Overlapped Execution

# "The Brick Wall"

**"Power Wall + Memory Wall + ILP Wall = Brick Wall"**

**The Landscape of Parallel Computing Research:
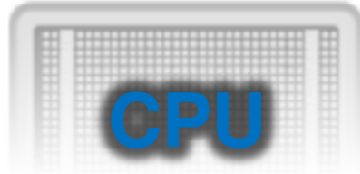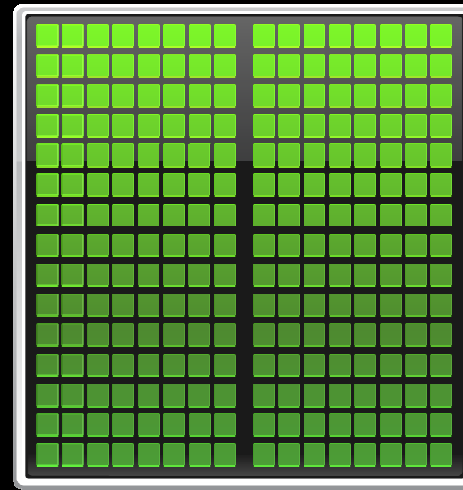A View from Berkeley**

**David Patterson et all.**

Performance (vs. VAX-)

52%/year

25%/year

20%/year

Source: Hennessy & Patterson, CAAQA, 4th Edition

10000
1000
100
10
1

1978  1980  1982  1984  1986  1988  1990  1992  1994  1996  1998  2000  2002  2004  2006

# Co-Processing
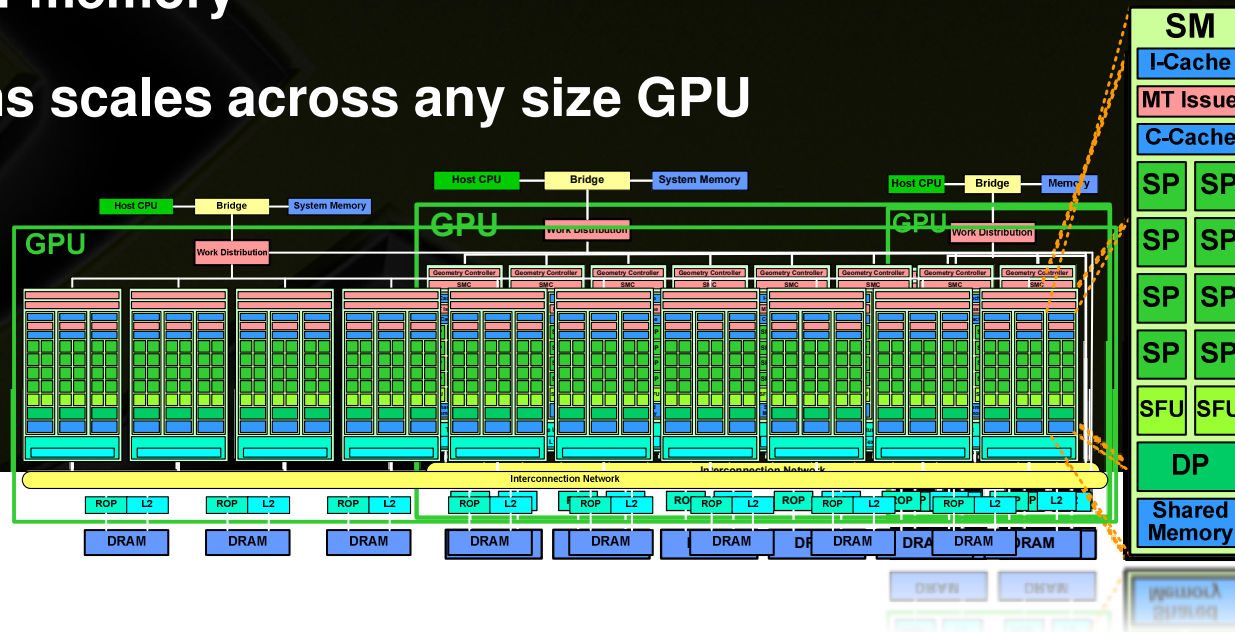## *The Right Processor for the Right Tasks*



CPU + GPU

# NVIDIA CUDA Parallel Computing Architecture

- **Many processors – eventually thousands**

- **Latency tolerant - execute 1000's of threads**

- **General load/store**

- **On-chip shared-memory**

- **CUDA programs scales across any size GPU**

# Rasterization & Ray Tracing

## Rasterization

- ### For each triangle
  - Find the pixels it covers
  - For each pixel: compare to closest triangle so far

Mapped to massively parallel GPU through DirectX or OpenGL

## Classical Ray Tracing

- ### For each pixel
  - Find the triangles that might be closest
  - For each triangle: compute distance to pixel

Mapped to massively parallel GPU through NVIDIA OptiX
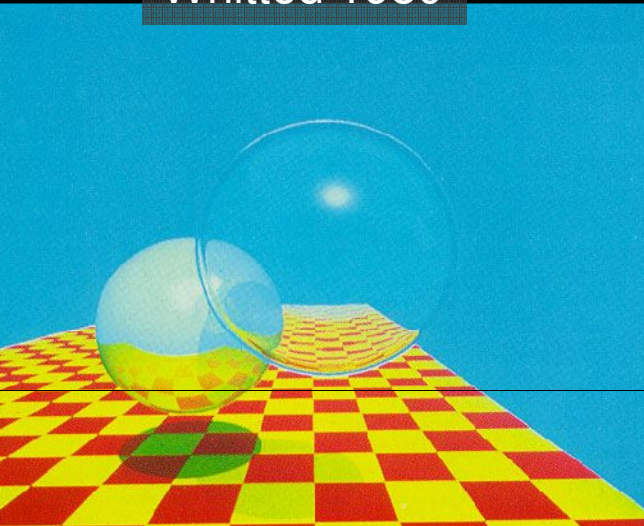
# Why ray tracing?

- **Ray tracing unifies rendering of visual phenomena**
  - fewer algorithms with fewer interactions between algorithms

- **Easier to combine advanced visual effects robustly**
  - soft shadows
  - subsurface scattering
  - indirect illumination
  - transparency
  - reflective & glossy surfaces
  - depth of field
  - ...
- **But: resource intensive, challenging to make fast**
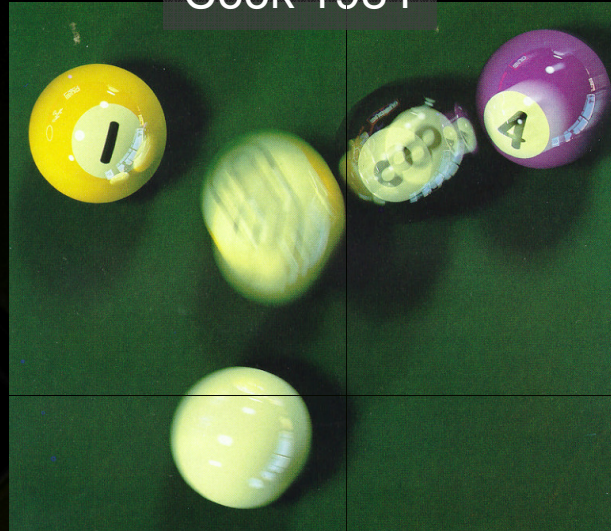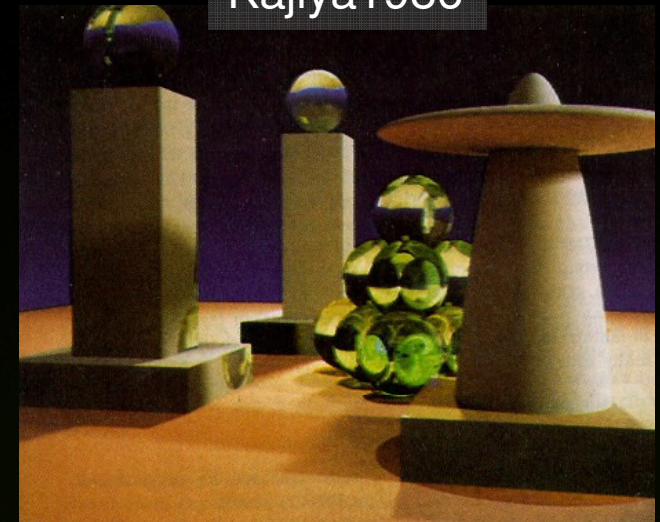
# Ray tracing regimes

Whitted 1980

Cook 1984

Kajiya1986

- Mirror reflections
- Perfect refractions
- Hard shadows

- 2-20 rays per pixel

- Depth of field
- Motion blur
- Soft shadows
- Glossy reflections
- 20-200 rays per pixel

- Indirect illumination
- Caustics
- Physical accuracy

- 200-$10^5$ rays per pixel

# OptiX Examples
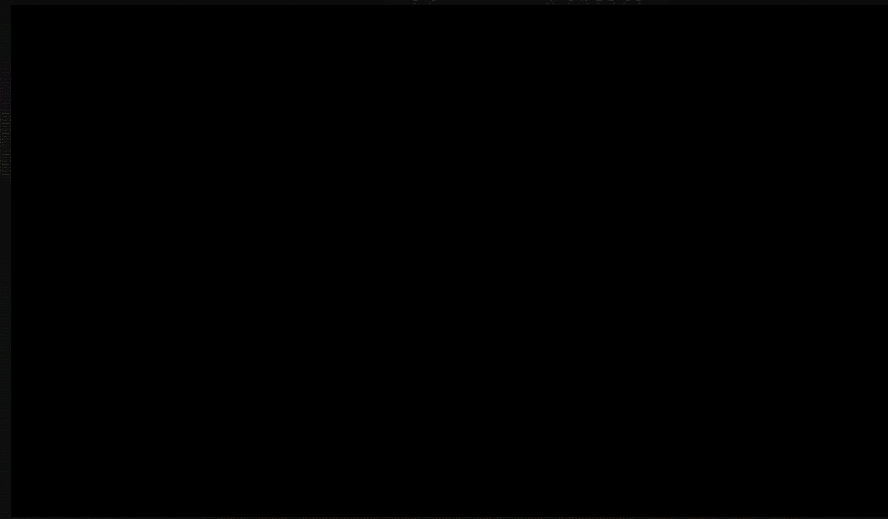


Interactive

# OptiX Examples



**Interactive**

**Progressive**

Cook 1984

•Mirror reflections     •Depth of field     •Indirect Illumination

•Perfect refractions     •Motion blur     •Caustics

•Hard shadows     •Soft shadows     •Physical accuracy

    •Glossy reflections

•1-10 rays per pixel     •10-100 rays per pixel     •100-10$^5$ rays per pixel
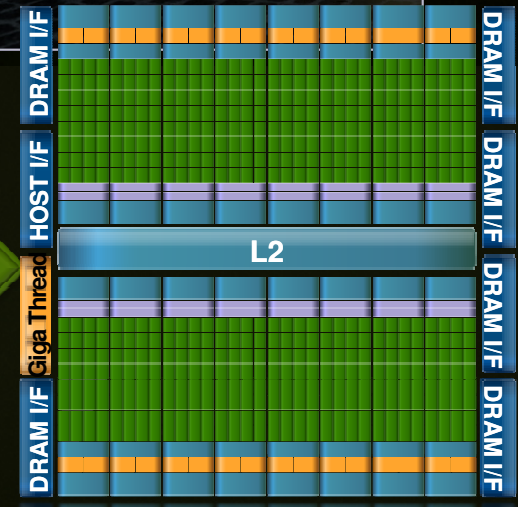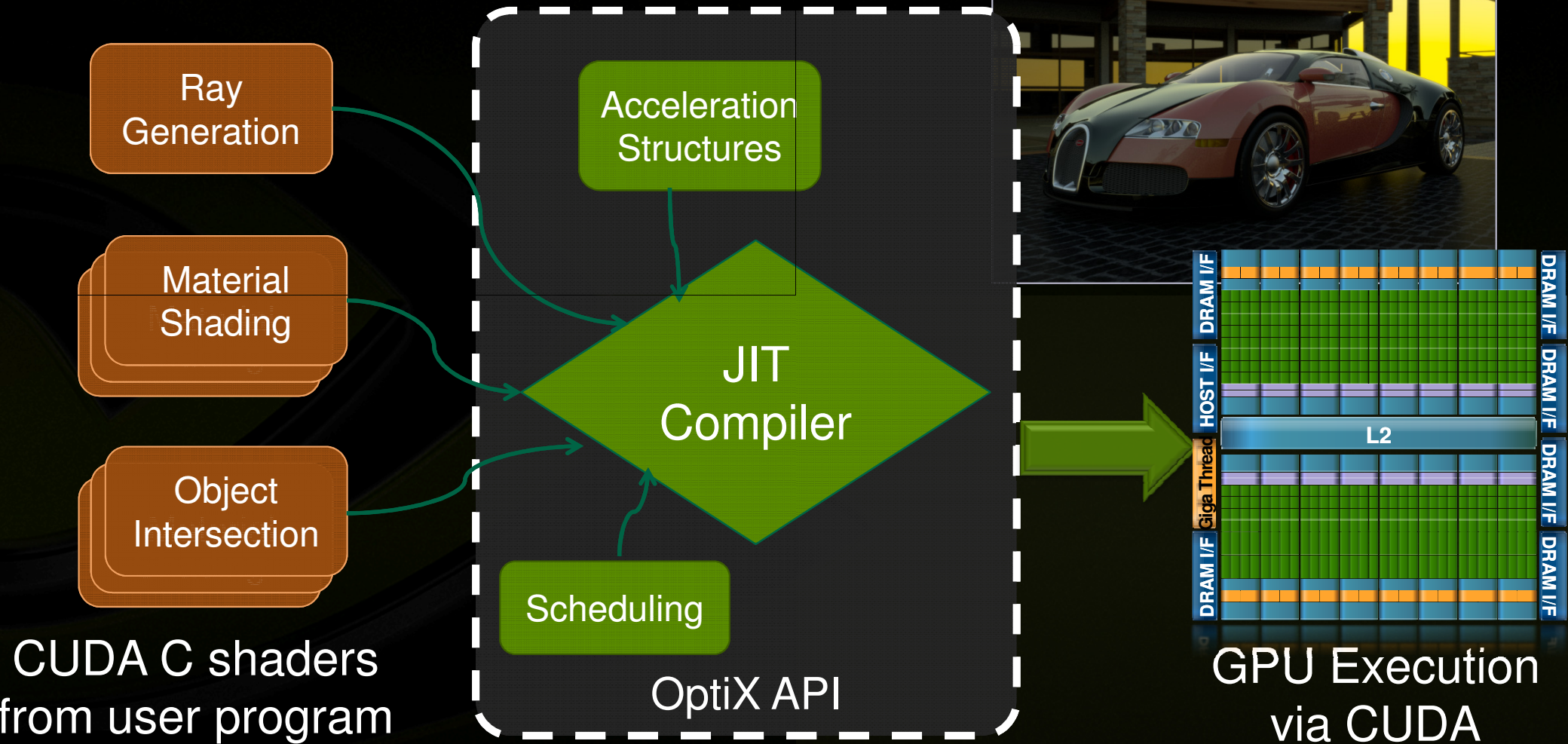
# Programmable Operations

## Rasterization

- **Fragment**

- **Vertex**
- **Geometry**

## Ray Tracing (OptiX)

- **Closest Hit**
- **Any Hit**
- **Intersection**
- **Selector**
- **Ray Generation**
- **Miss**
- **Exception**

The ensemble of programs defines the rendering algorithm

(or collision detection algorithm, or sound propagation algorithm, etc.)
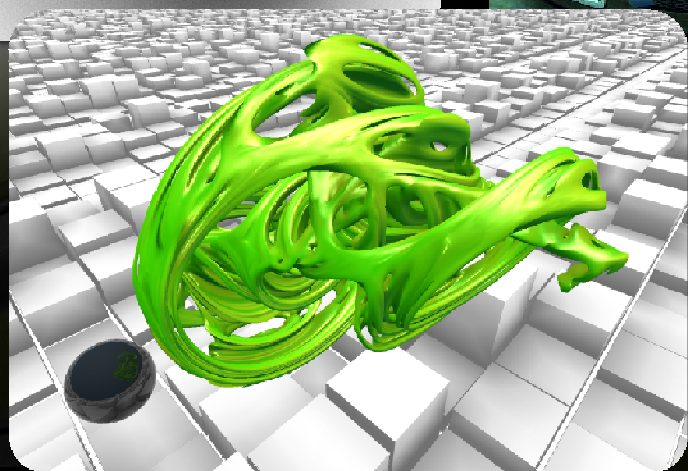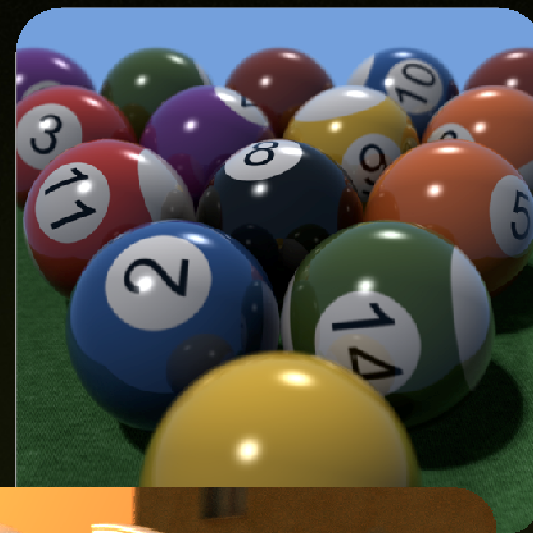
# OptiX Functional Overview
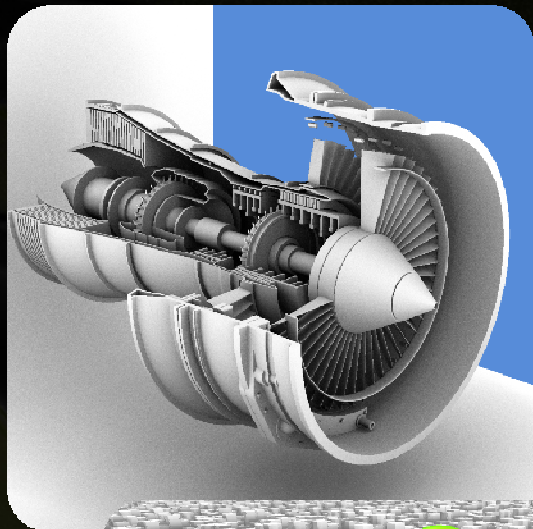
**Ray Generation**

**Material Shading**

**Object Intersection**

CUDA C shaders
from user program

Acceleration Structures

JIT Compiler

Scheduling

OptiX API

GPU Execution
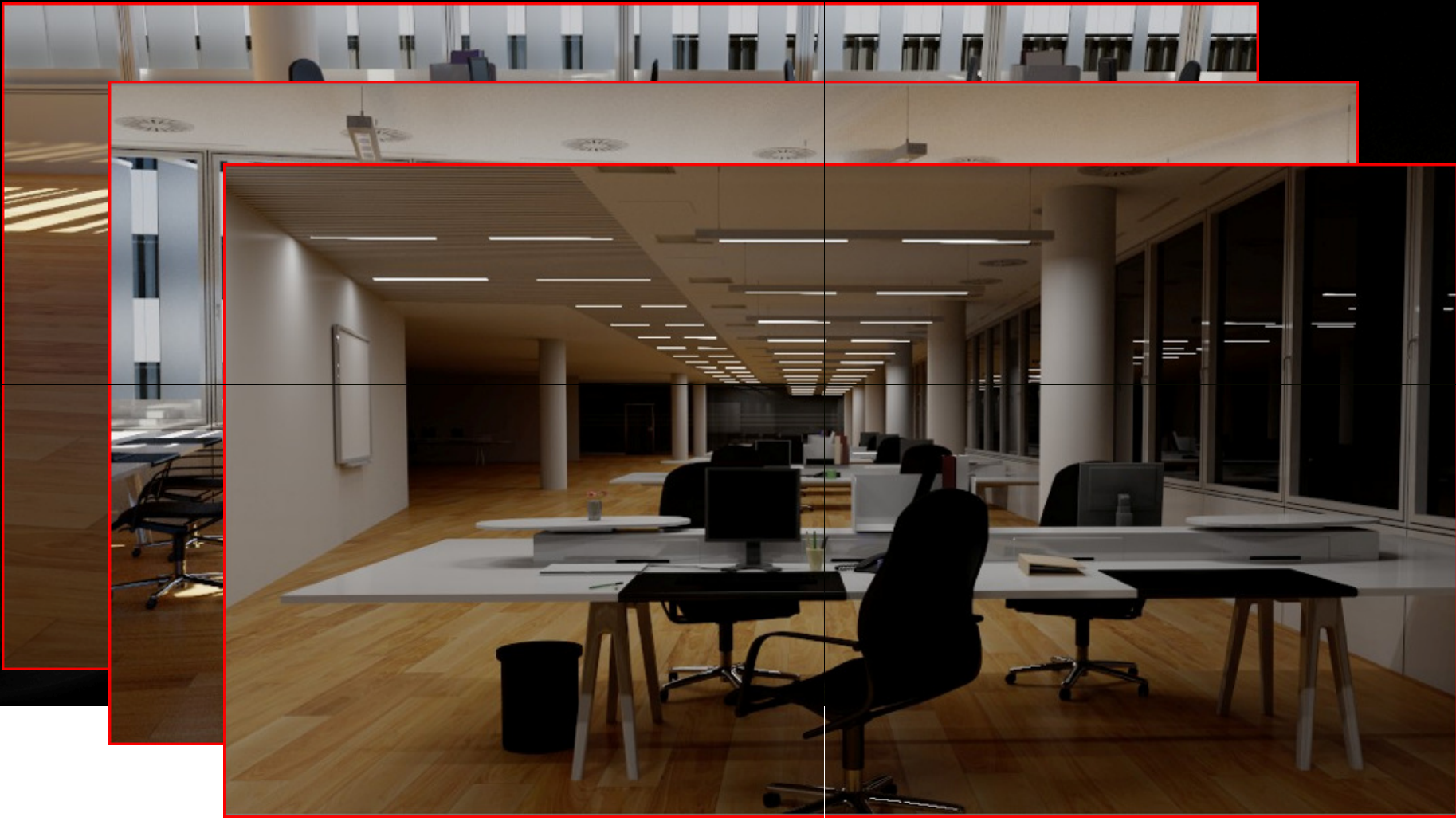via CUDA

# OptiX SDK

**Available now at http://www.nvidia.com**

# iray® Light Transport Simulation

- No more renderer introduced artifacts

- No more parameter tweaking

- Deterministic quasi-Monte Carlo integro-approximation

  - consistent is faster than unbiased

  - converges unconditionally

  - exactly reproducible independent of parallelization

- High precision ray tracing

- Scales across architectures

  - takes optimal advantage of GPU compute power
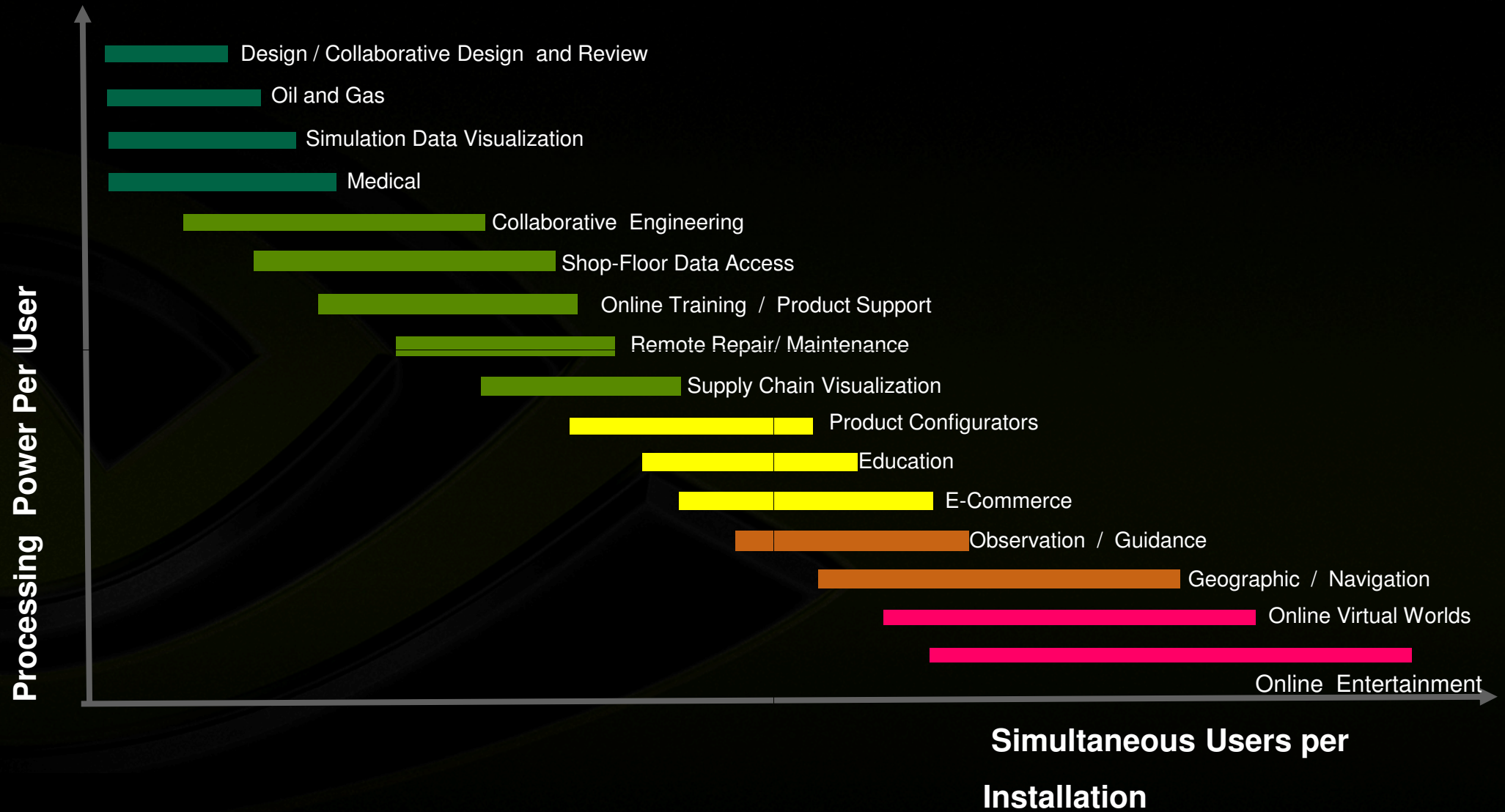
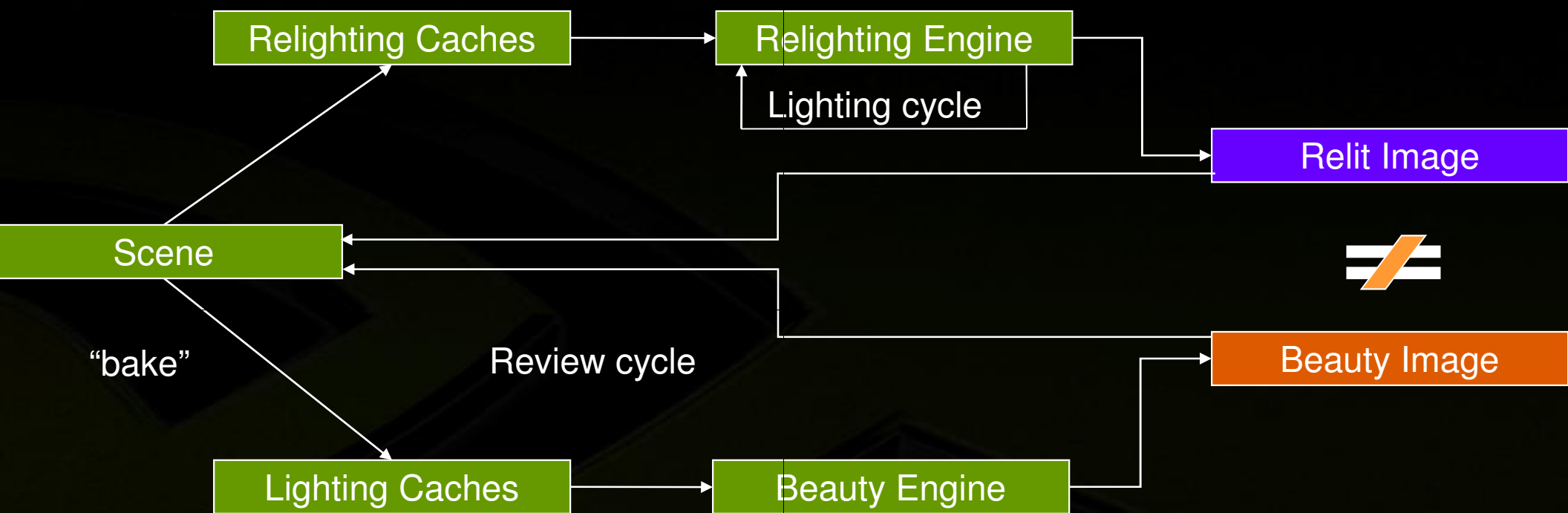# iray® Light Transport Simulation

# RealityServer

**Software Platform for 3D Applications and Web Services**

- **remote interaction regardless of 3D data complexity**

- **thin clients including mobile devices**

- **collaboration**

- **data security**

- **built on top of Distributed Computing Environment (DiCE)**
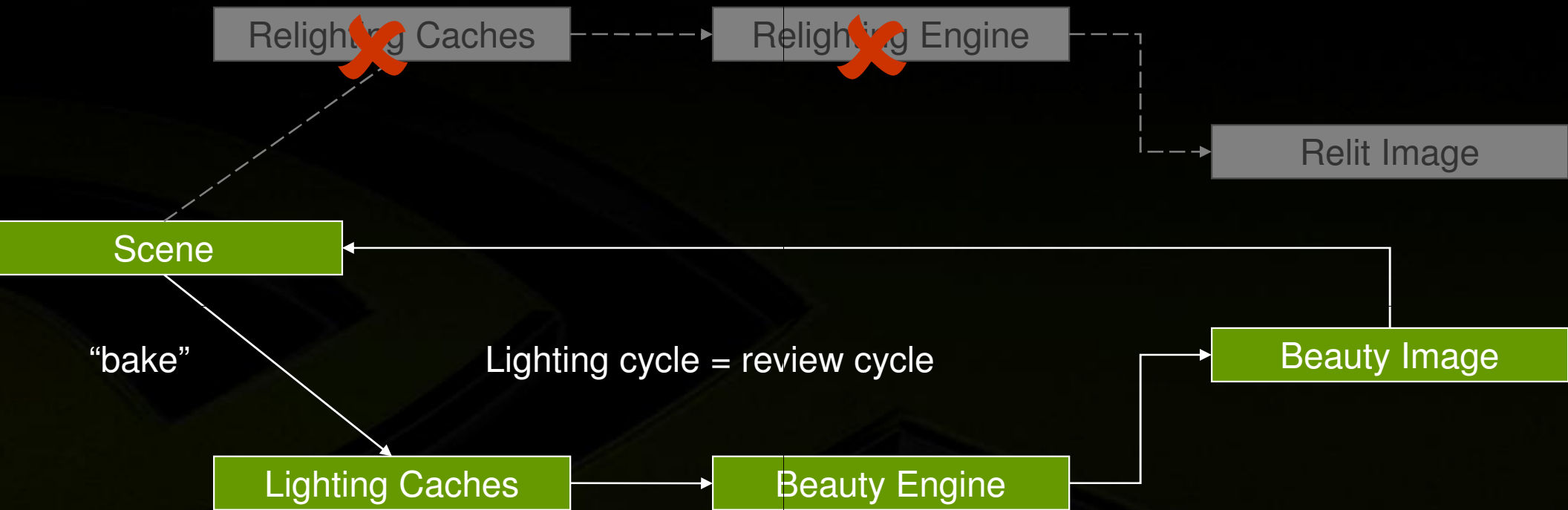
- **includes iray® GPU simulation technology**

# RealityServer Applications

**Processing Power Per User**

- Design / Collaborative Design and Review
- Oil and Gas
- Simulation Data Visualization
- Medical
- Collaborative Engineering
- Shop-Floor Data Access
- Online Training / Product Support
- Remote Repair / Maintenance
- Supply Chain Visualization
- Product Configurators
- Education
- E-Commerce
- Observation / Guidance
- Geographic / Navigation
- Online Virtual Worlds
- Online Entertainment

**Simultaneous Users per Installation**

# A Typical Lighting Cycle



Relighting Caches → Relighting Engine → Relit Image

Lighting cycle

Scene ≠

"bake"

Review cycle → Beauty Image

Lighting Caches → Beauty Engine

# Beauty lighting

Relighting Caches  ❌ - - - - → Relighting Engine  ❌ - - - - →

Relit Image

Scene

"bake"

Lighting cycle = review cycle

Beauty Image

Lighting Caches → Beauty Engine

# PantaRay – Precomputation Engine

- **Special purpose caching engine based on data structures optimized for GPU deployment**

- **Extreme scene complexity**
  - **Normal shots in the 10-100M polygons**
  - **Complex shots over 1G polygons**

- **Accelerate beauty pass with raytraced PRT caches**

- **Compared to the previous method, the tractability limit is increased by 2 orders of magnitude (100x)**

Voxel Rendering: Sibenik

# Voxel Rendering: Why Voxels?

- **Future challenges**
  - **More geometric detail**
  - **Unique assets**
  - **Authoring, capturing**
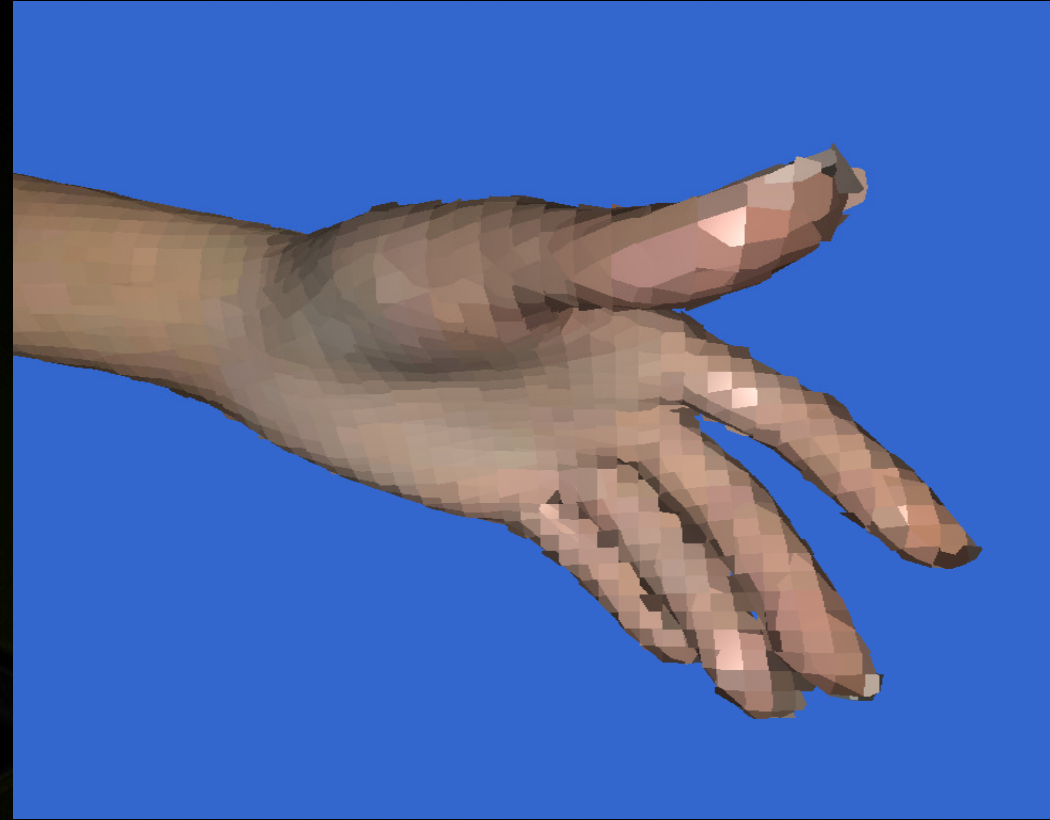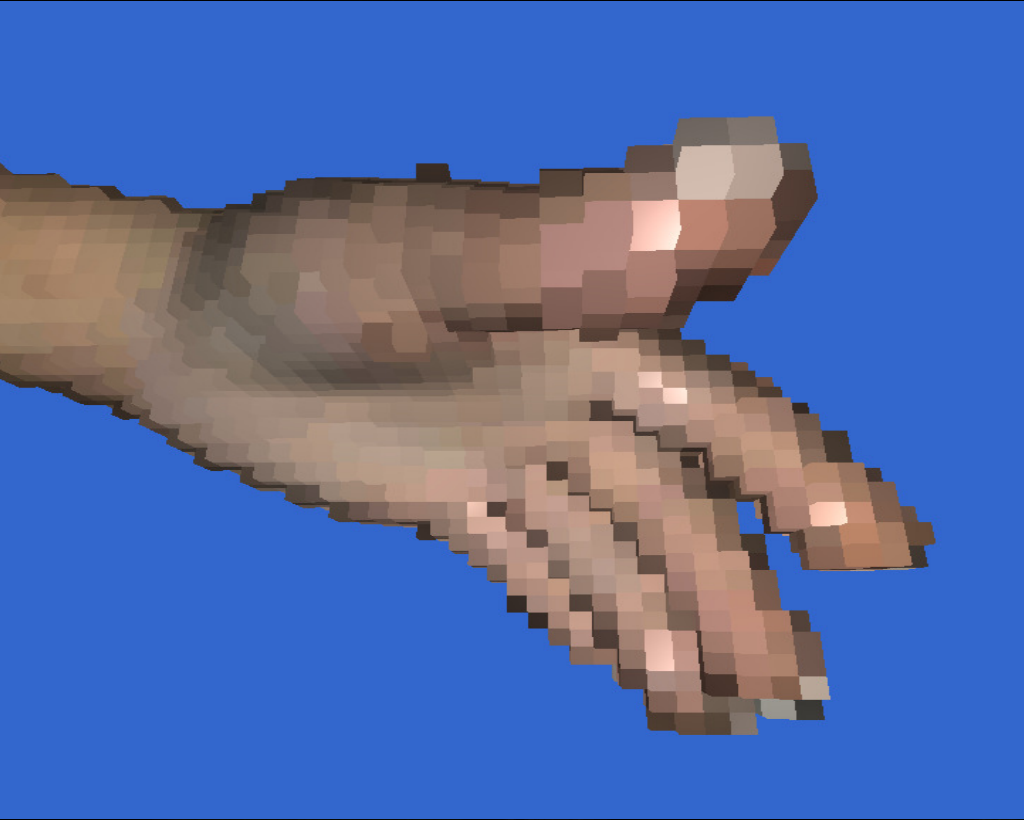  - **Ray tracing**
- **Hot Idea?**
  - **Leading game developers looking at voxels**
  - **Web press, game technology fans talking it up**
- **Historical trends**
  - **Raw data wins**
  - **Simplicity turns into efficiency – think of Z-buffer!**

# Contours



ote: low-resolution voxelization

# Post-Process Blur



No blur

With blur
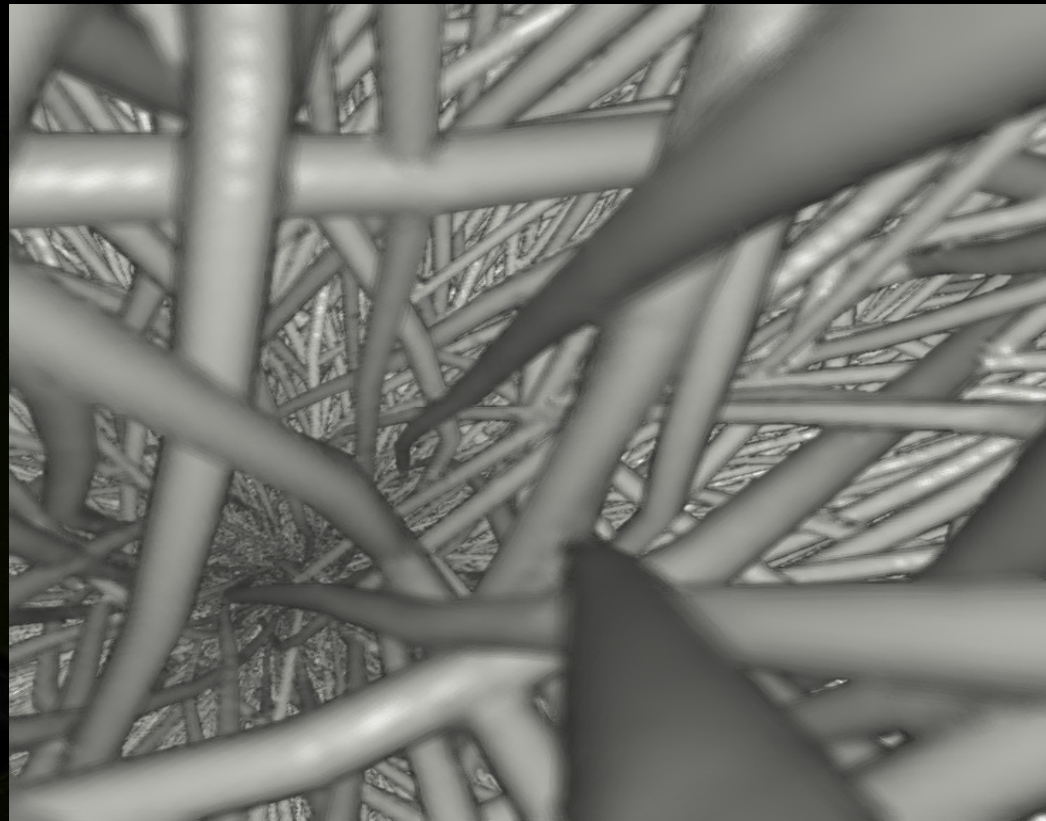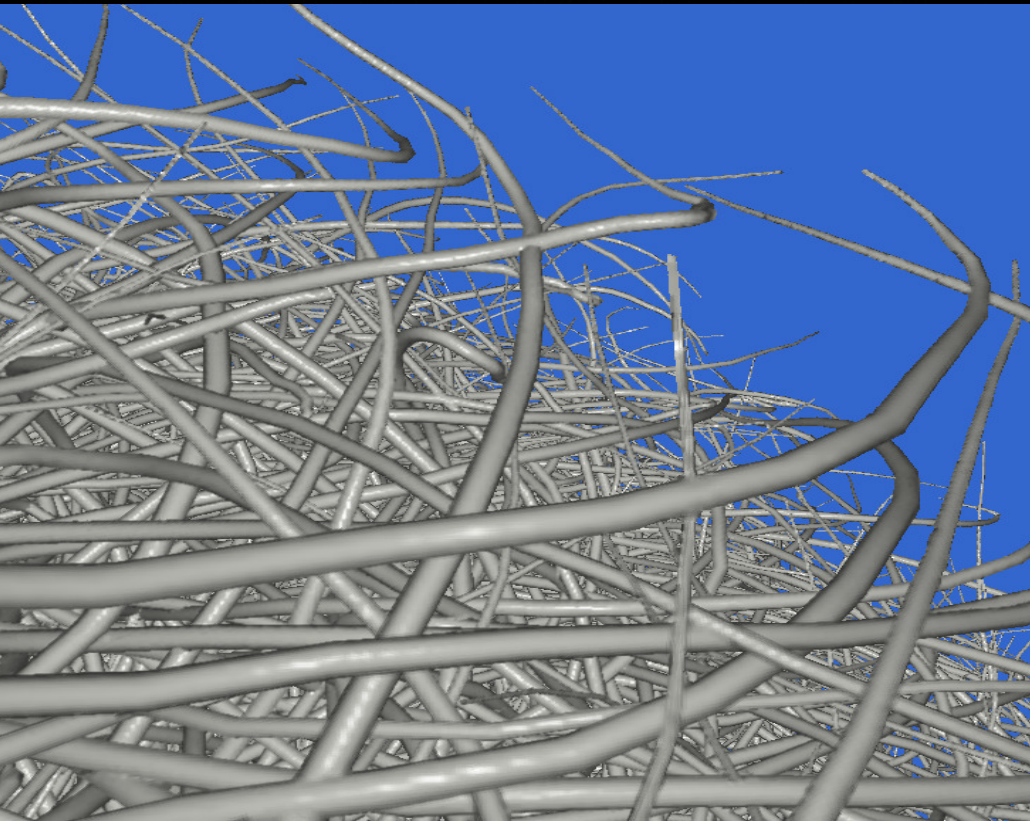
ote: low-resolution voxelization
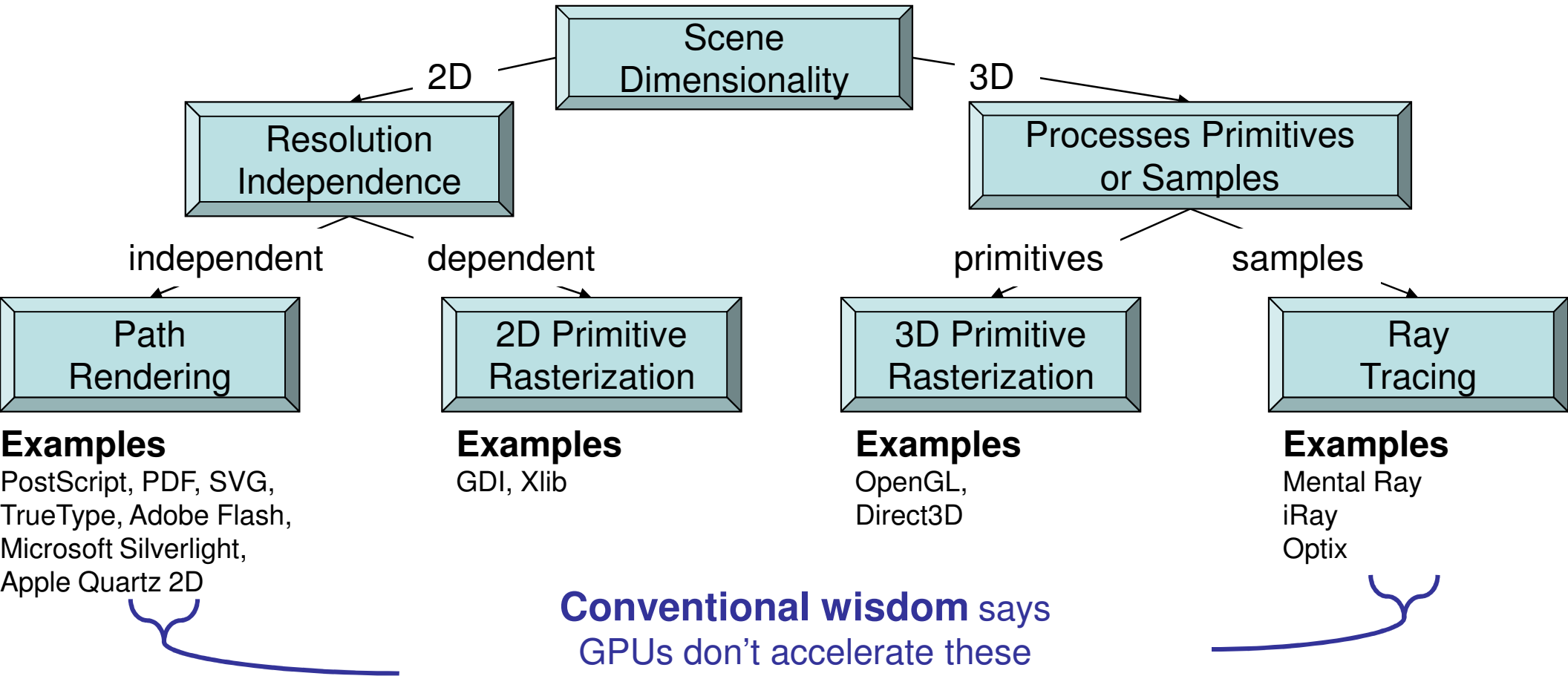
Hairball, 11 levels: 10:28, 1552 MB, 28.4 Mrays/s

# Resolution Highlights



Hairball, 11 levels: 1552 MB

# Taxonomy of Rendering

```
                        ┌──────────────────┐
                        │      Scene       │
                   2D   │  Dimensionality  │   3D
            ┌───────────┤                  ├───────────┐
            ▼           └──────────────────┘           ▼
┌──────────────────┐                      ┌──────────────────────┐
│    Resolution    │                      │ Processes Primitives │
│   Independence   │                      │      or Samples      │
└──────────────────┘                      └──────────────────────┘
```

| independent | dependent | primitives | samples |
|---|---|---|---|
| **Path Rendering** | **2D Primitive Rasterization** | **3D Primitive Rasterization** | **Ray Tracing** |

**Examples**
PostScript, PDF, SVG,
TrueType, Adobe Flash,
Microsoft Silverlight,
Apple Quartz 2D

**Examples**
GDI, Xlib

**Examples**
OpenGL,
Direct3D

**Examples**
Mental Ray
iRay
Optix

**Conventional wisdom** says
GPUs don't accelerate these
rendering paradigms
*As usual, the conventional wisdom is WRONG* ☺

# Path Rendering Standards

**Document Printing and Exchange**

PDF

Adobe PostScript

*Open XML Paper (XPS)*

**Resolution-Independent Fonts**

*OpenType*

*TrueType*

**Immersive Web Experience**

*Flash*

Microsoft Silverlight

*Scalable Vector Graphics*

HTML 5

**2D Graphics Programming Interfaces**

*Java 2D API*

*QtGui API*

Quartz 2D GRAPHICS

*Mac OS X 2D API*

OpenVG

*Khronos API*

**Office Productivity Applications**

Microsoft Office

OpenOffice.org

*Adobe Illustrator*

*Open Source Inkscape*

# GPU vs. CPU
# Path Rendering Performance

## Frames/second rendering Celtic dogs



**"Celtic round dogs" scene**
1 very complex path
5,031 cubic Bezier commands
29,068 coordinates

|  | GeForce 8600 GT x4 | GeForce 8600 GT x8 | GeForce 8600 GT x16 | Cairo software renderer | Qt software renderer |
|---|---|---|---|---|---|
| 1000x1000 pixel window | 250 | 126 | 93 | 13.1 | 1.1 |
| 500x500 pixel window | 725 | 350 | 272 | 22.2 | 2.7 |

# GPU-Accelerated Path Rendering Exploits GPU Advantages

- Use stencil to track filled and stroked path coverage
  - GPUs rasterize stencil only at 2x to 3x faster than normal peak shaded rendering rate
  - [Loop & Blinn 2005] discard shader techniques support rasterizing paths with curved segments
  - No tessellation required
    - Avoids expensive, sequential process
- Then shade path coverage with conventional GPU-accelerated shading
  - GPUs far better at shading than CPUs
- Antialiasing - Hardware multisampling very efficient

*All path content below rendered by GPU without tessellation*

# Even Font Rendering by GPU

- GPU-accelerated font rendering means the end of glyphs rasterized to bitmaps
  - GPU can render directly from outlines

hese glyph samples are GPU rendered from their PostScript and TrueType outlines

# GPU Advantage: Example of Clipping to an Arbitrary Path

- Path rendering standards require path clipping
  - A rendered path can be clipped to another arbitrary path
- Expensive on the CPU
  - Must compute intersection of two arbitrary paths
- Cheap on the GPU
  - Stencil buffering is very efficient

Dragon masked by letter Psi

The computational graphics continuum



**"Just" programmable shading: DX, OGL**

The computational graphics continuum



**"Just" programmable
shading: DX, OGL**

**"Pure" compute-based
graphics: CUDA, OptiX**

The continuum "Beyond Programmable Shading"



Interesting middle ground!

"Just" programmable shading: DX, OGL

"Pure" compute-based graphics: CUDA, OptiX

**Image Space Photon Mapping**



"ust" programmable
ading: DX, OGL

"Pure" compute-based
graphics: CUDA, OptiX

"Image-Space Photon Mapping",
Morgan McGuire, David Luebke. High Performance Graphics 2009

**DPL16**    I sped up the wipes to 0.2s, see what you think.
David Luebke, 7/31/2009

# Photon Mapping Time (seconds)

e.g., [Jensen 01, Purcell et al. 03, Zhou et al. 08]

**1ˢᵗ Bounce**  **2ⁿᵈ**  **3ʳᵈ**  **4ᵗʰ...**  Build Tree  **Last Bounce** ("Estimate Radiance")

# Image Space Photon Mapping (*milliseconds*)

GPU Bounce Map  CPU Trace  GPU Photon Volumes

*Data Transfer*

**DPL17**      technically its rebalancing the tree, right? doesn't traditional PM build the tree incrementally then balance it to speed up the knn-search?
David Luebke, 7/31/2009

**DPL18**      I rearranged the arrows and animation a bit to emphasize the single point of projection for first and last bounces, and to make the "appear" animations into ultra-fast (0.2s) fades.
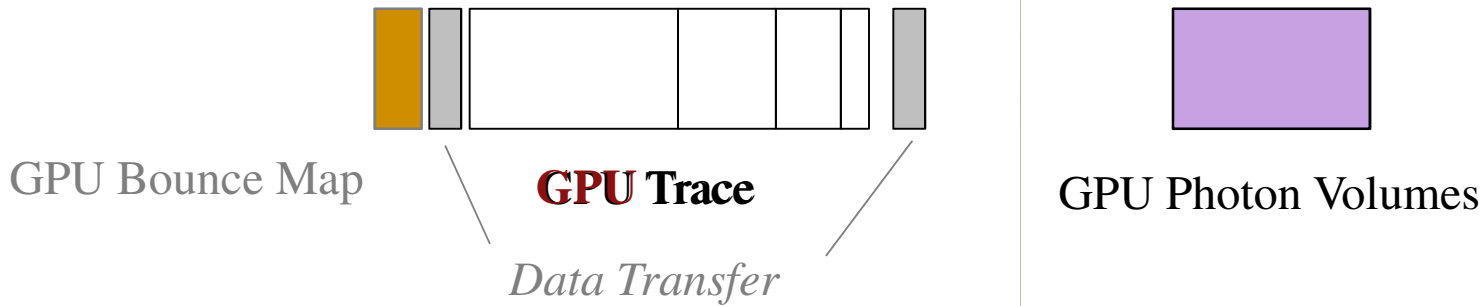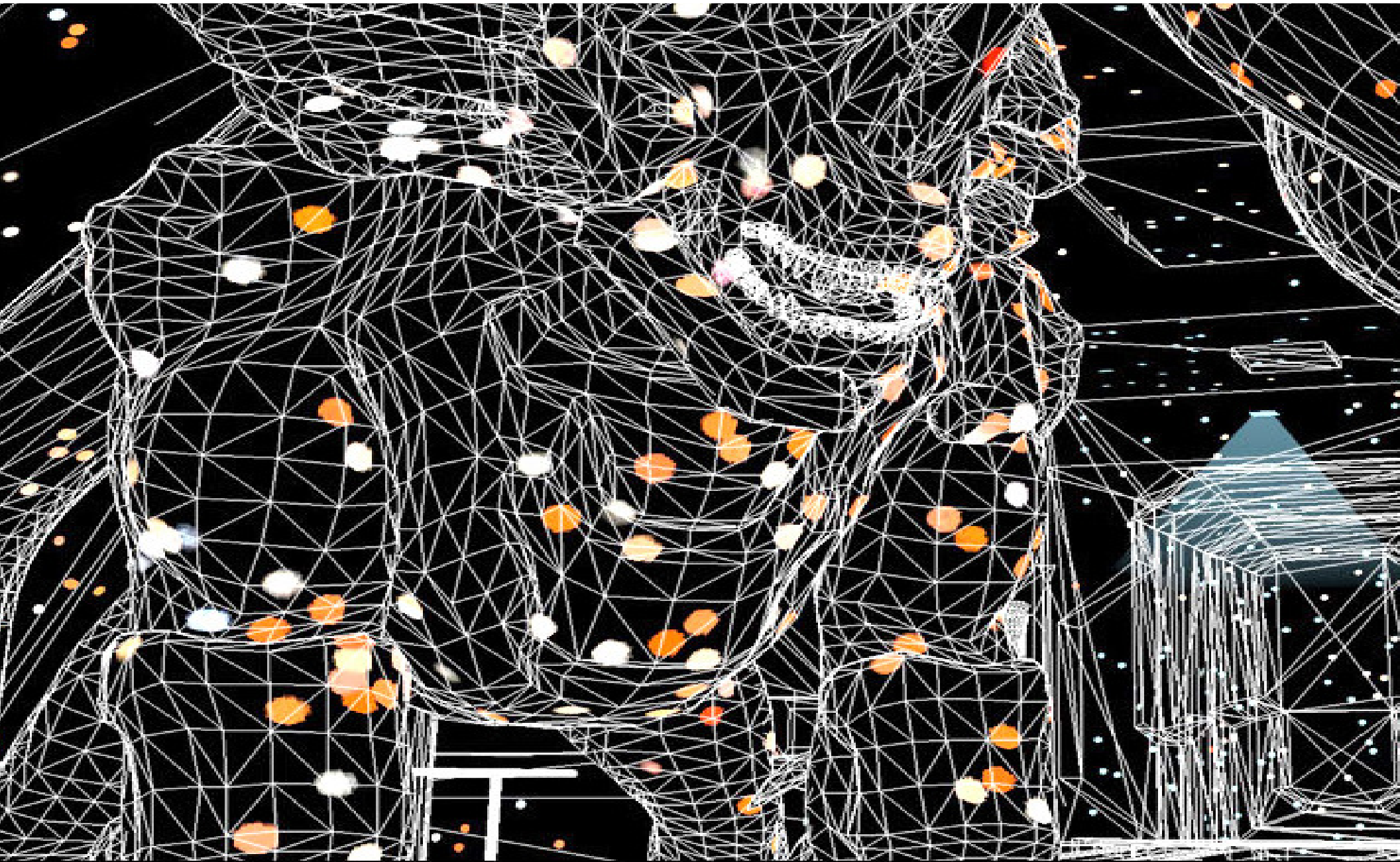David Luebke, 7/31/2009

# Photon Mapping Time (seconds)

e.g., [Jensen 01, Purcell et al. 03, Zhou et al. 08]

**1st Bounce**   **2nd**   **3rd**   **4th...**   Build Tree   **Last Bounce**
**("Estimate Radiance")**

# OptiX-enabled Image Space Photon Mapping

GPU Bounce Map   **GPU Trace**   GPU Photon Volumes

*Data Transfer*

**DPL19**	technically its rebalancing the tree, right? doesn't traditional PM build the tree incrementally then balance it to speed up the knn-search?
David Luebke, 7/31/2009

**DPL20**	I rearranged the arrows and animation a bit to emphasize the single point of projection for first and last bounces, and to make the "appear" animations into ultra-fast (0.2s) fades.
David Luebke, 7/31/2009

Direct

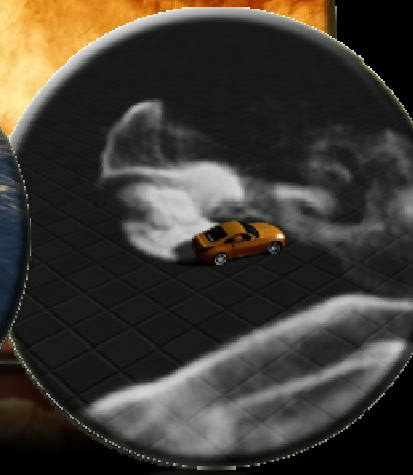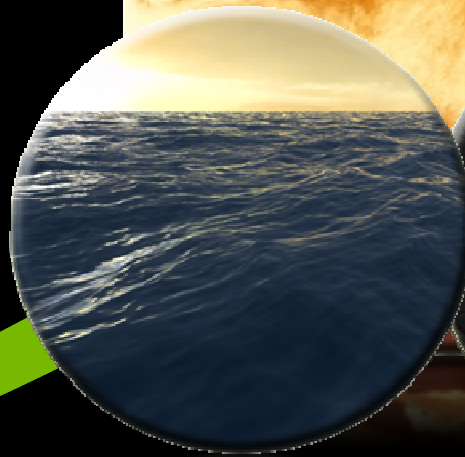Indirect

Direct +
Indirect

Direct +
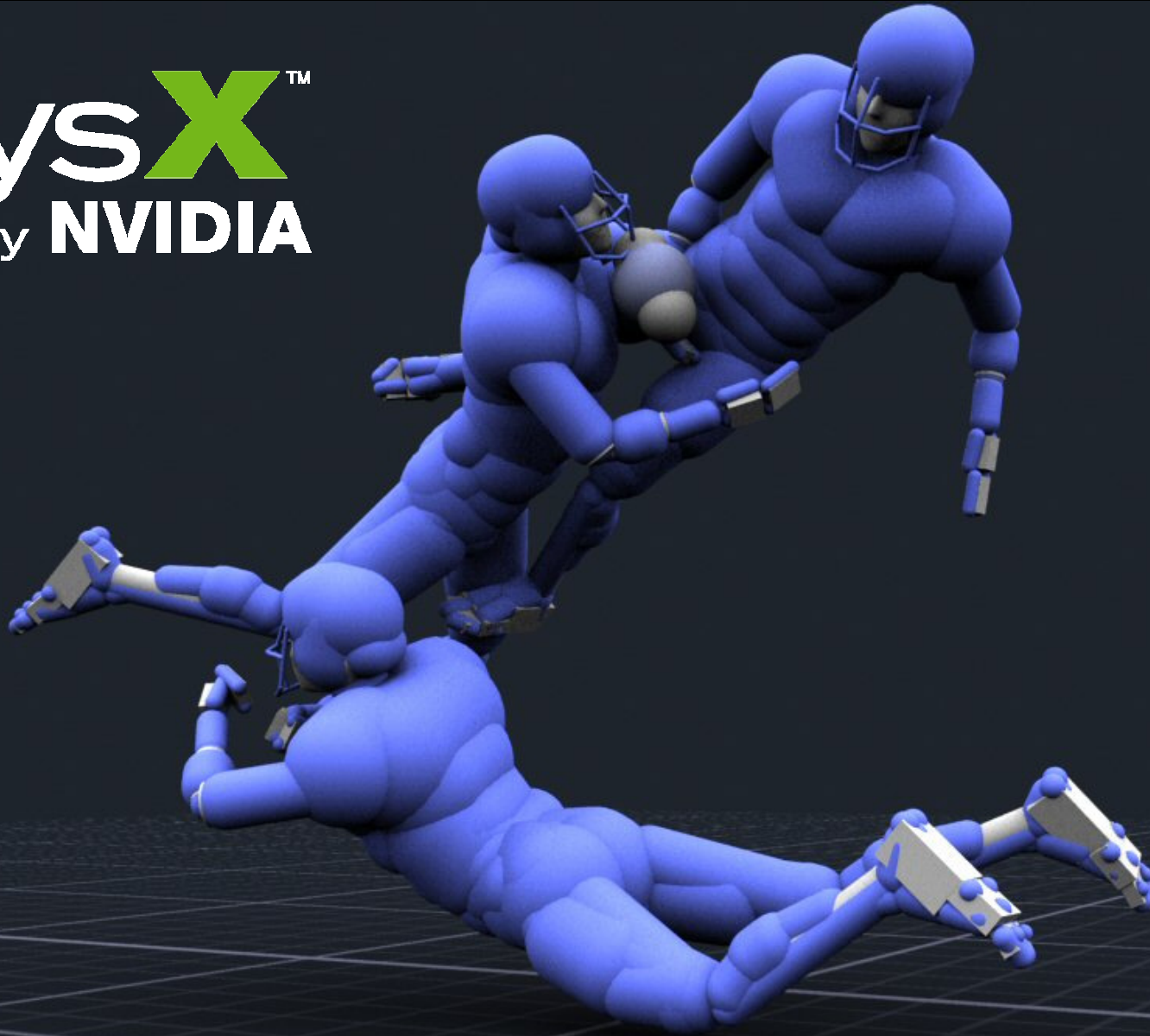Ambient

Direct +
Indirect

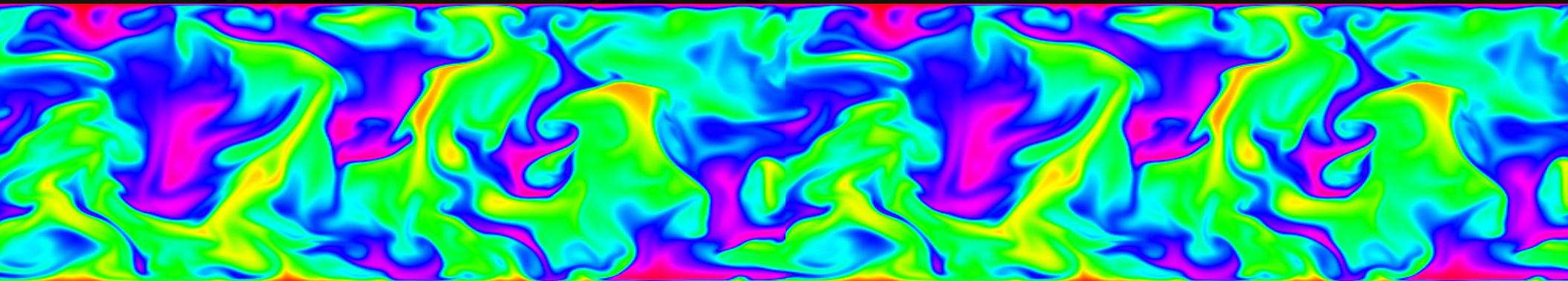# The Next Big Thing – Physics
*Simulate Amazing Worlds*

# Real-time fluid effects

**Complex fluid-drive motion is all around**

- Car exhaust, dust storms, rolling mist, steam, smoke, fire, contrails, bubbles in water, …

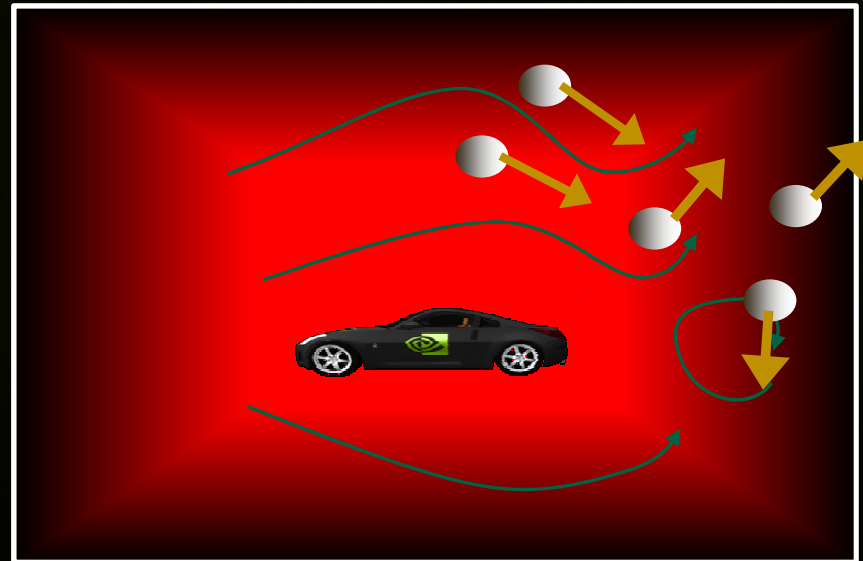**Goal: Add this level of realism to games**

**Problem: Turbulent motion is computationally intensive!**

# Solution: GPUs are computational monsters!

**Calculate near-field fluid on grid**
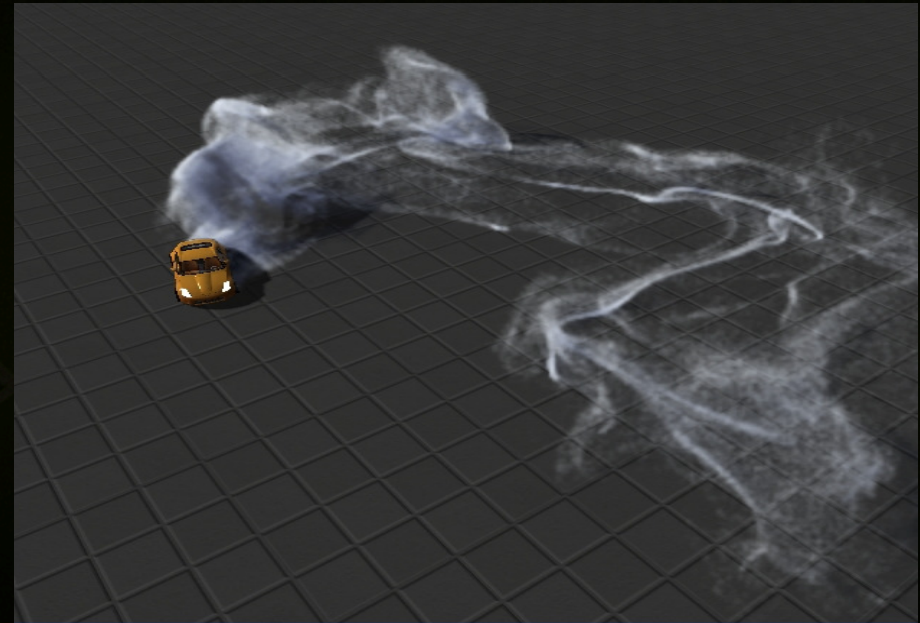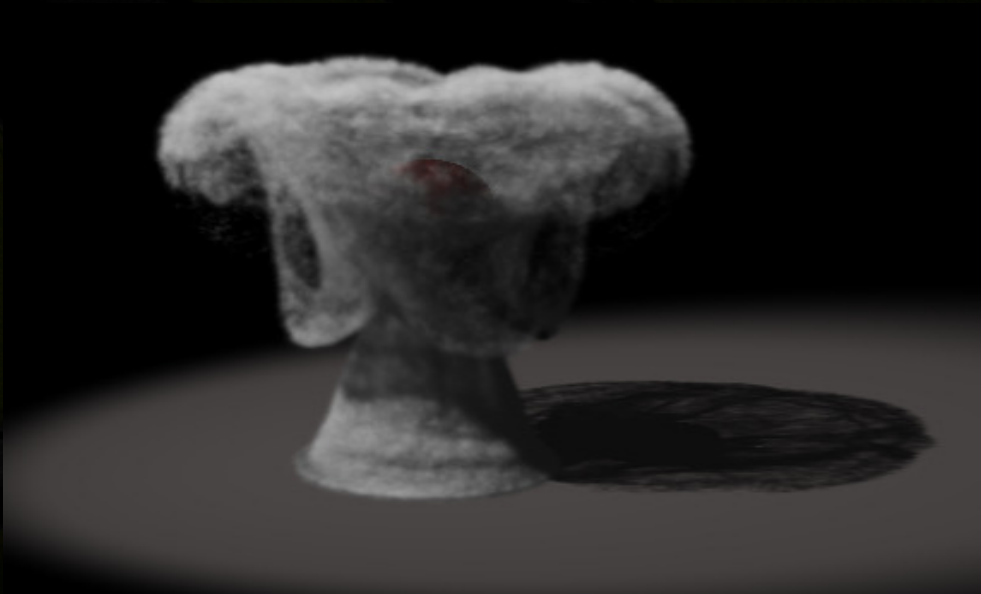**Fluid velocities drive particle motions**

1. Calculate Fluid Velocities on Regular Grid
   2nd-Order Accurate CUDA Multigrid Solver

2. Interpolate Fluid Velocities onto Particles
   3D Interpolation in CUDA

3. Advance Particles
   CUDA Particle System

4. Render Particles
   CUDA - OpenGL Interop

# APEX Turbulence

**Interactive CFD Solution + Volume Rendering**

Interactive Fluid-Particle Simulation
Using Translating Eulerian Grids.

# Co-Processing
*The Right Processors for the Right Tasks*

**2015 Projection**

| | | |
|---|---|---|
| CPU-Alone | 1.2^6 | 3X |
| CPU+GPU | 50 * 1.5^6 | 570X |

*Heterogeneous Parallel Processing will (continue to) Fundamentally Change the way we do Graphics*